

Análise de séries temporais

por similaridade e alinhamento não
linear com Dynamic Time Warping



Busca por Similaridade com DTW

Provendo eficiência

Eficiência computacional

O termo “eficiência” é bastante genérico e pode se referir a tempo de execução, memória, etc

Na Computação, usamos a análise assintótica como ferramenta

- Mas vamos pular detalhes

DE vs DTW

Vamos considerar duas séries temporais de mesmo comprimento n .

A distância euclidiana faz uma operação para cada par de observação

$$ed(x, y) = \sqrt{\sum_{k=1}^n (x_k - y_k)^2} \sim 3 \cdot n \text{ operações}$$

Portanto, o número de operações realizadas é proporcional ao tamanho da entrada (do algoritmo), que no caso é n .

Dizemos, portanto, que a ED é $O(n)$ ¹ ou que o algoritmo é linear.

↳ big-oh

1. Tem muita simplificação aqui, mas aceitamos assim, mesmo na Computação

DE vs DTW

A DTW precisa preencher a matriz de custo, que é $n \times n$. Ou seja, a DTW precisa fazer $O(n^2)$ operações.

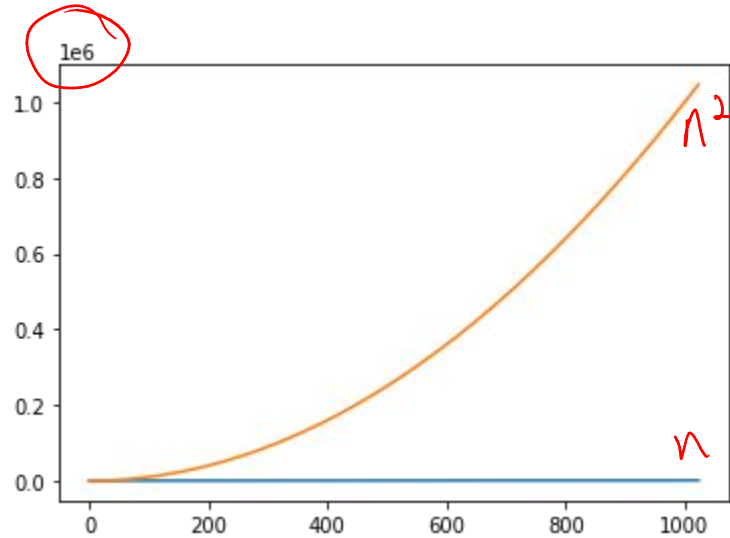
- Portanto, dizemos que a DTW é quadrática (nesse caso, tanto em tempo de execução quanto em memória)
- Usando as bandas de Sakoe-Chiba de tamanho w , ficamos com $O(nw)$

$$\begin{aligned} i &= 1 \dots n \\ j &= i-w \dots i+w \end{aligned} \quad O(n \cdot w)$$

Isso é muito importante, pois o tempo de execução depende do número de operações.

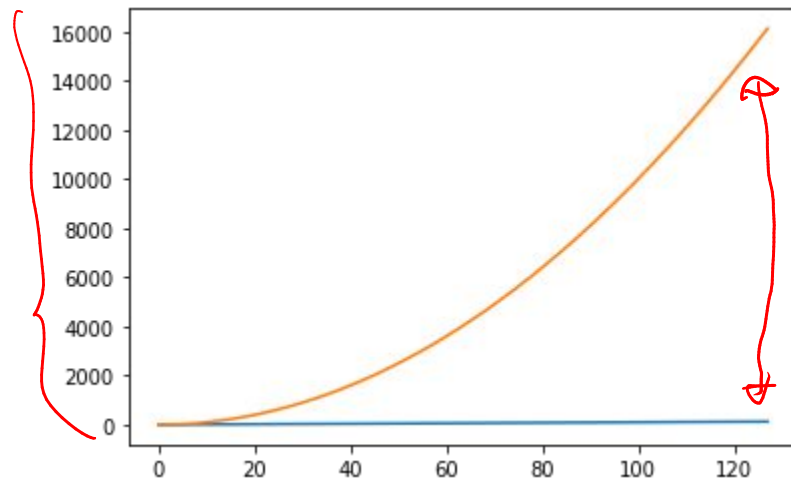
Linear vs Quadrático

$n = 1024$



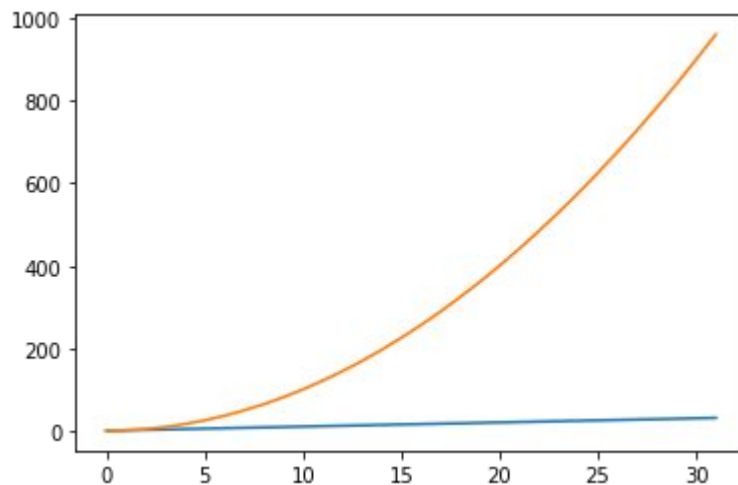
Linear vs Quadrático

$n = 128$



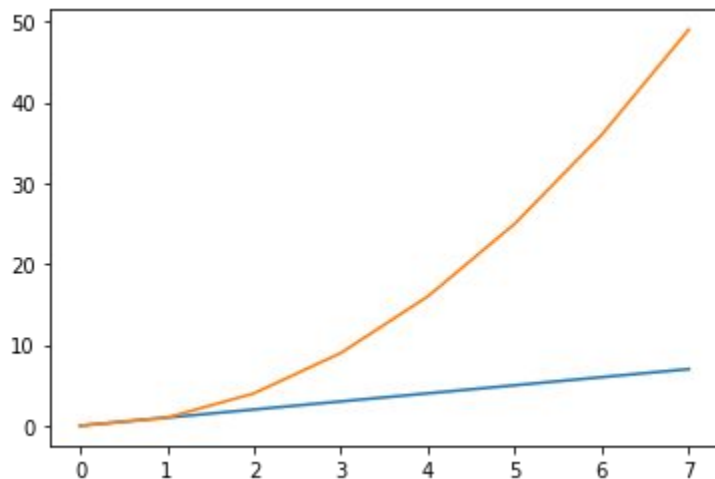
Linear vs Quadrático

$n = 32$



Linear vs Quadrático

$n = 8$



DE vs DTW

Para comparar um único par de séries temporais, a diferença não é significativa, na prática.

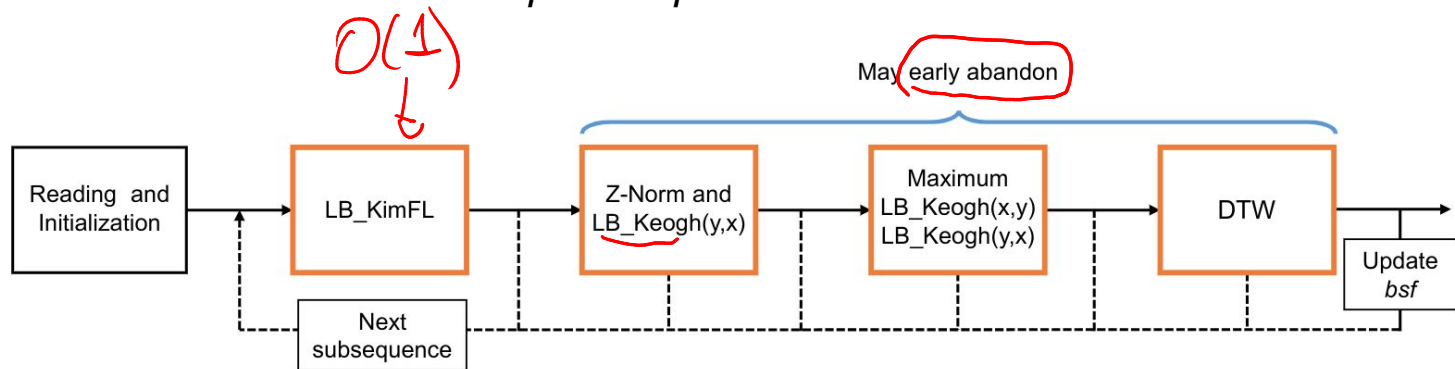
Mas e se eu for comparar vários, como na busca por similaridade?

- O valor n é o tamanho da consulta
- Multiplica-se um fator l , que é o comprimento da série de referência
 - e l costuma ser bem grande (esse é exatamente o objetivo)

$$O(l \cdot \text{complex}(\text{dist}))$$

UCR Suite

Portanto, para **escalar** uma busca por similaridade baseada em DTW, é preciso de várias técnicas de *speed-up*.



Um ponto importante antes de dissecar o pipeline: *best-so-far*

UCR Suite

O primeiro ponto importante é não ~~tirar~~ ^{calcular} a raiz da DTW

- Sim, eu dei uma mentida até aqui
- ED e DTW tem uma raiz quadrada e eu omiti a do DTW

Motivo...

UCR Suite

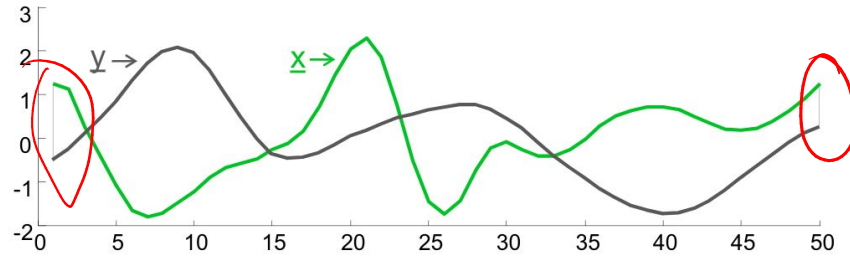
LB-Kim-FL

iterações 500
 $d(bsf) = 7,45$

$bsf = \text{best-so-far}$

$O(1)$

$LB_Kim_FL(500) = 8,49$



Um *lower bound* é uma medida garantidamente menor que a *DTW distance*.

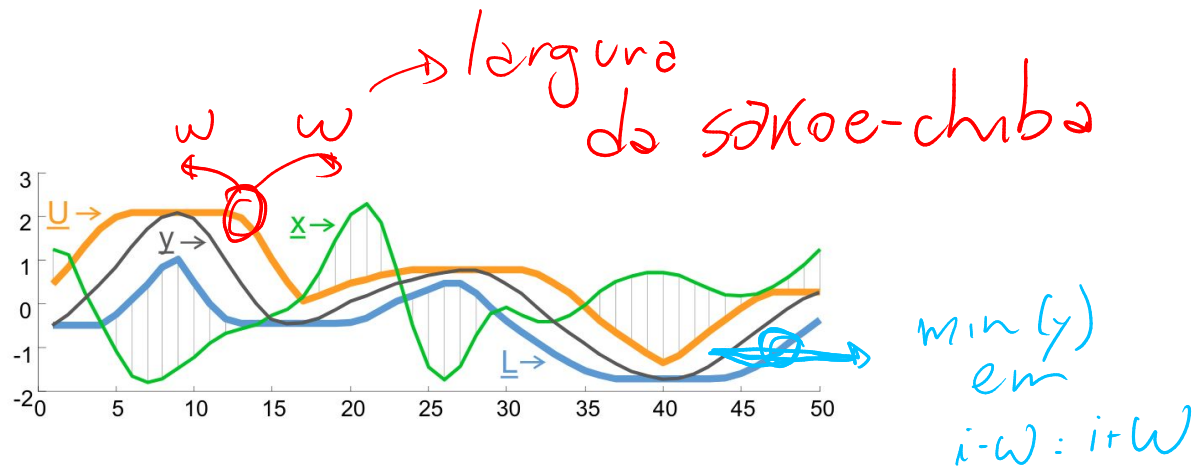
- Para ser útil, precisa ser rápido de ser calculado

$$LB_Kim_FL = (x_1 - y_1)^2 + (x_n - y_n)^2 \leq DTW(X, Y)$$

Considerando-se
a indexação
começando em 1

UCR Suite

LB-Keogh

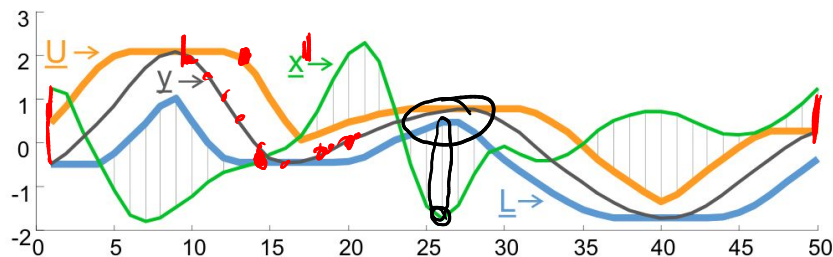


A *suite* considera LB_Keogh(X,Y) e LB_Keogh(Y,X)

UCR Suite

LB-Keogh

$$\sum \left\{ \begin{array}{l} \text{if } L_i \leq x_i \leq U_i \\ \min \left\{ c(x_i, U_i), c(x_i, L_i) \right\}, \text{ c.c.} \end{array} \right.$$

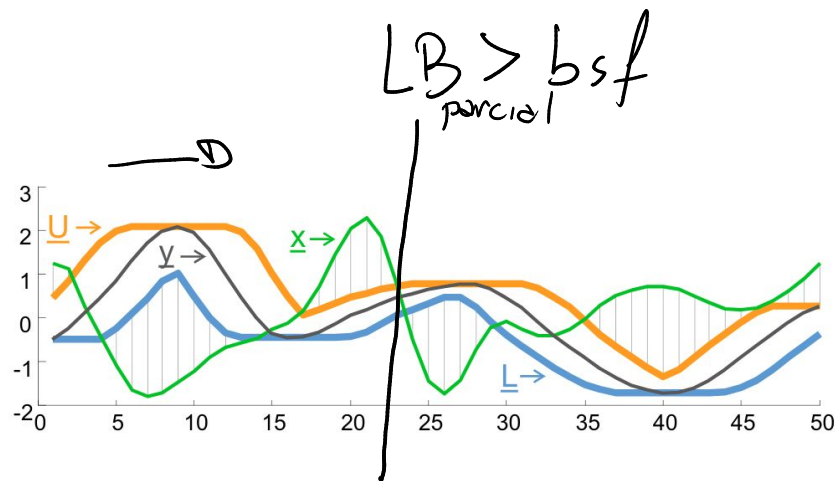


A suite considera $LB_Keogh(X,Y)$ e $LB_Keogh(Y,X)$

$LB_Keogh > LB_Kim$

UCR Suite

LB-Keogh



A *suite* considera $LB_Keogh(X,Y)$ e $LB_Keogh(Y,X)$

UCR Suite

Early abandon do lower bound

- Ao calcular o LB-Keogh, é possível abandonar precocemente
- Se o cálculo parcial ultrapassa o *bsf*, podemos descartar a subsequência

Early abandon da z-normalização

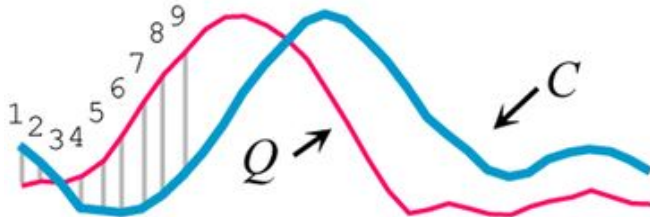
- Considera-se a busca sempre z-normalizada por subsequência
- Em vez de calcular de início, alterna-se com o LB

+ early abandon
de DTW

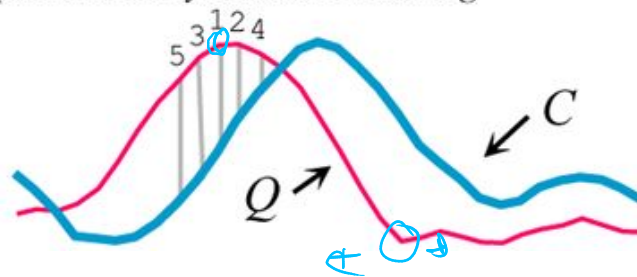
UCR Suite

Re-ordenar o cálculo

Standard early abandon ordering



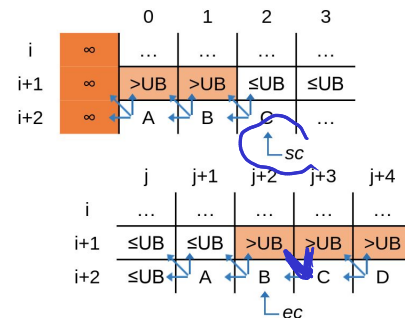
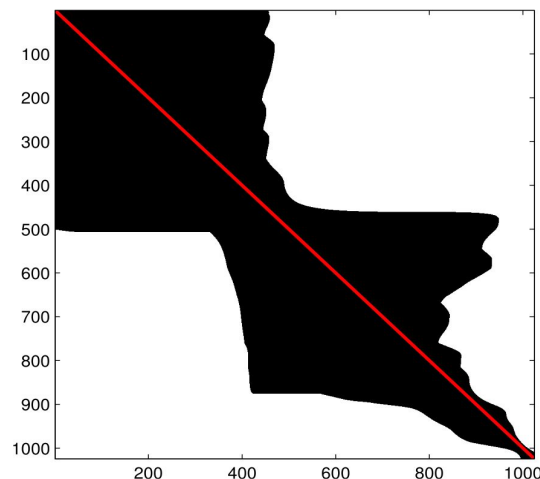
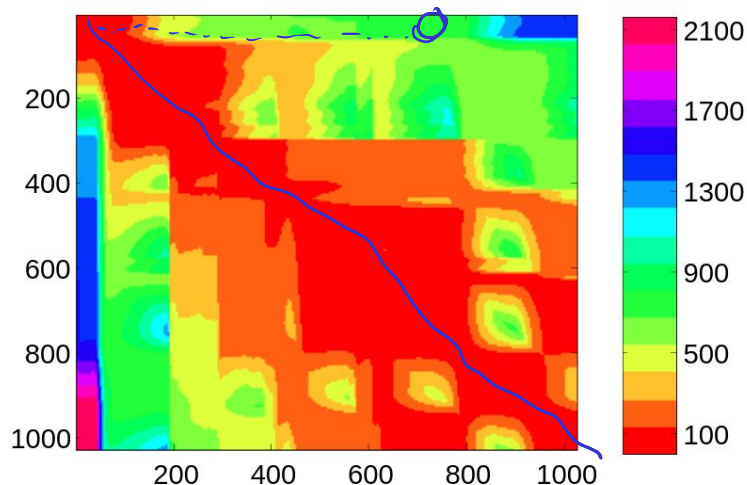
Optimized early abandon ordering



UCR-USP Suite

Mais otimizações

UB = dist. Euclidean
 Grande busca: $UB = bsf$
 Pruned DTW



UCR-USP Suite

Vamos implementar isso tudo? Não, claro. Mas vamos à prática